

AU NO. DOC FILE COPY



CENTER FOR CYBERNETIC STUDIES

D D C

DE STERNING

JUL 28 1978

DISCUSSIVES

EVEN

TO STERNING

TO ST

The University of Texas Austin, Texas 78712

This document has been approved for public release and sale; its distribution is unlimited.



79 07 24 005

LEVEI S 9 1 Research Report CCS-318 AD A 056 A SPECIAL PURPOSE LINEAR PROGRAMMING ALGORITHM FOR OBTAINING LEAST ABSOLUTE VALUE ESTIMATORS IN A LINEAR MODEL WITH DUMMY VARIABLES. by Ronald D. Armstrong Edward Frome** *The University of Texas at Austin, Austin, TX 78712 **Oak Ridge Associated Universities, Oak Ridge, TN 37830 NOD014-75-C-0569 NSF-MCS707-00100 This research was partly supported by NSF Grant MCS707-00100, and ONR Contract N00014-75-C-0569 with the Center for Cybernetic Studies, The

University of Texas. Reproduction in whole or in part is permitted for any purpose of the United States Government.

CENTER FOR CYBERNETIC STUDIES

A. Charnes, Director Business-Economics Building, 203E The University of Texas Austin, Texas 78712 (512) 471-1821

This document has been approved for public release and sale; is distribution is unlimited.

406 197

A SPECIAL PURPOSE LINEAR PROGRAMMING ALGORITHM FOR OBTAINING LEAST ABSOLUTE VALUE ESTIMATING IN A LINEAR MODEL WITH DUMMY VARIABLES

Ronald D. Armstrong

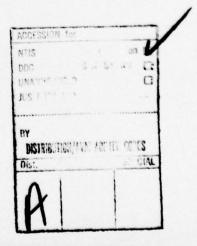
The University of Texas Austin, Texas

Edward Frome

Oak Ridge Associated Universities Oak Ridge, Tennessee

ABSTRACT

Dummy (0,1) variables are frequently used in statistical modeling to represent the effect of certain extraneous factors. This paper presents a special purpose linear programming algorithm for obtaining least-absolute-value estimators in a linear model with dummy variables. The algorithm employs a compact basis inverse procedure and incorporates the advanced basis exchange techniques available in specialized algorithms for the general linear least-absolute-value problem. Computational results with a computer code version of the algorithm are given.



INTRODUCTION

The standard linear regression model is usually written as follows:

 $y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + ... + x_{im}\beta_m + \varepsilon_i$, i = 1, ..., n, (1) where $x_i = (x_{i1}, x_{i2}, ..., x_{im})$ are known values of the independent variables, y_i is an observed value of a dependent random variable, ε_i is an error term and the β_i 's are the unknown parameters.

It is often desirable to include parameters in the model that represent the "effect" of different levels of one or more "factors". This can be done using dummy (0,1) variables to represent the levels of the factor, and is usually referred to as the analysis of covariance - e.g., see Searle (1971, chap. 4). Throughout this paper we assume that each observation is affected by at most one of the K levels of a single factor. The regression model then becomes

$$y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + ... + x_{im}\beta_m + d_{i1}\alpha_1 + d_{i2}\alpha_2 + ... + d_{ik}\alpha_k + \epsilon_i$$
 (2) where

$$d_{ik} = \begin{cases} 1 \text{ if kth level of factor is present,} \\ 0 \text{ otherwise, k=1, ..., K,} \end{cases}$$
 (3)

and the α_k 's are unknown parameters representing the effect of each level of the factor on the response.

The classical procedure for estimating the unknown parameters is to solve the following least squares problem.

Minimize
$$\sum_{j=1}^{n} (y_j - \sum_{j=1}^{m} x_{ij}\beta_j - \sum_{k=1}^{K} d_{ik}\alpha_k)^2.$$
 (4)

Computational techniques for solving (4) can be found in Searle (1971, chap. 4). Strong theoretical justification (Graybill (1961)) can be made for the least squares estimates when certain assumptions on the distributions of the random variables are made. However, in the presence of fat-tailed distributions or outliers, least absolute value (LAV) estimates may be recommended. Empirical studies (see Barrodale (1968), Glahe and Hunt (1970), Kiountouzis (1973)) comparing least squares and LAV estimators have demonstrated

the worth of the LAV criterion. The LAV problem may be stated:

Minimize
$$\sum_{j=1}^{n} |y_{j} - \sum_{j=1}^{m} x_{ij} \beta_{j} - \sum_{k=1}^{K} d_{ik} \alpha_{k}|.$$
 (5)

It is generally recognized that the LAV problem may be solved efficiently with a special purpose linear programming algorithm (see Barrodale and Roberts (1973)). The main purpose of this paper is to develop a further refinement of the linear programming approach to handle dummy variables. The advantages of the refinement are a reduction in computer processing time, and a reduction in computer storage.

Section 2 reviews the LAV algorithm of Barrodale and Roberts (1973) and develops the theory necessary to take advantage of the structure arising in the LAV problem when dummy variables are present. Section 3 discusses the computer code implementation and presents some computational results. Section 4 presents sample problems where LAV estimation is used.

2. THE LINEAR PROGRAMMING ALGORITHM

2.1 Review of the Barrodale and Roberts' Algorithm

We will begin this section with a brief review of a revised simplex implementation of the Barrodale and Roberts' algorithm.

A more detailed description of this procedure - generalized to handle linear constraints - is given by Armstrong and Hultz (1977).

Charnes, Cooper and Ferguson (1955) demonstrate how the LAV problem associated with (1) can be written as the following linear programming (LP) problem.

Minimize
$$z = \sum_{i=1}^{n} (P_i + N_i)$$

subject to

$$X\beta + IP - IN = Y$$

P > 0, N > 0

where X is the n by m matrix (x_{ij}) , Y is the vector of y_i 's, P and N are, respectively, the vectors of positive and negative deviations.

A basis for this LP problem is formed by m independent rows of

X. We denote this submatrix by X_B and its inverse by X_B^{-1} . We define IB to be the ordered set of row indices forming X_B (the row index of X corresponding to the first row of X_B is the first element of IB and so on) and NB is the index set of <u>nonbasic</u> rows.

We make the transformation $\lambda = X_B \beta$ and the problem becomes:

Minimize
$$z = \sum_{i=1}^{n} (P_i + N_i)$$

subject to

$$\lambda_{q} + P_{i(q)} - N_{i(q)} = y_{i(q)}, q = 1, 2, ..., m$$

 $X_{i}X_{B}^{-1}\lambda + P_{i} - N_{i} = y_{i}, i \in NB$

where i(q) is the q-th element of IB.

A current basic solution consists of

$$\overline{\lambda}_{\mathbf{q}} = y_{\mathbf{i}(\mathbf{q})}, \ \mathbf{q} = 1, 2, \dots, m$$

$$\overline{\beta} = X_{\mathbf{B}}^{-1} \overline{\lambda}$$

$$P_{\mathbf{i}(\mathbf{q})} = N_{\mathbf{i}(\mathbf{q})} = 0, \ \mathbf{q} = 1, 2, \dots, m$$

$$\overline{P}_{\mathbf{i}} = \begin{cases} y_{\mathbf{i}} - X_{\mathbf{i}} \overline{\beta}; \ \mathbf{i} \in NB, \ y_{\mathbf{i}} - X_{\mathbf{i}} \overline{\beta} > 0 \\ 0; \ \mathbf{i} \in NB, \ y_{\mathbf{i}} - X_{\mathbf{i}} \overline{\beta} < 0 \end{cases}$$

$$\overline{N}_{\mathbf{i}} = \begin{cases} X_{\mathbf{i}} \overline{\beta} - y_{\mathbf{i}}; \ \mathbf{i} \in NB, \ X_{\mathbf{i}} \overline{\beta} - y_{\mathbf{i}} > 0 \\ 0; \ \mathbf{i} \in NB, \ X_{\mathbf{i}} \overline{\beta} - y_{\mathbf{i}} \leq 0 \end{cases}$$

To assist in obtaining the LP <u>reduced costs</u> or rate of change achieved by removing a row from the basis, we define

$$\sigma_i = \text{sgn}(y_i - X_i \overline{\beta}), i \in NB$$

 $\sigma_i = +1 \text{ or } -1, \text{ when } y_i - X_i \overline{\beta} = 0, i \in NB.$

When the deviation for a nonbasic row is zero, an initial assignment of +1 or -1 to σ_i is arbitrary and thereafter the algorithm will determine the value.

The LP dual variables (see Wagner (1959) for a statement of the dual problem) associated with the basic constraints are given by

$$\pi = (\pi_1, \pi_2, ..., \pi_m) = (\sum_{i \in NB} \sigma_i X_i) X_B^{-1}.$$

The dual variables associated with the nonbasic constraints are +1 or -1. Thus, the optimality condition for the specialized LP algorithm of Barrodale and Roberts is

$$-1 \le \pi_i \le +1$$
, i = 1, 2, ..., m.

Let us assume that the optimality condition is not satisfied. Then, there exists a component of $\pi(\text{say}, \pi_r)$ with $|\pi_r| > 1$. The row $X_{i(r)}$ has now been labeled as leaving the basis. The value of $\sigma = \text{sgn}(\pi_r)$ indicates whether λ_r is to be increased or decreased. Because the objective is to minimize the sum of the residuals, λ_r should be increased if π_r is negative and λ_r should be decreased if π_r is positive. The objective value may not strictly decrease at each iteration when degeneracy (i.e., $\overline{P}_i = \overline{N}_i = 0$, i ϵ NB) is present. Degeneracy never seems to cause cycling in practice and can be resolved with the perturbation technique of Charnes (1952). We will ignore the implications of degeneracy for the remainder of this paper.

Once the algorithm has specified that $X_{i(r)}$ is to leave the basis, it must determine the row of X to enter the basis in the r-th position. This is accomplished through a partial sort of the ratios:

$$(y_i - X_i\overline{\beta})/(\rho X_i X_B(r)), \ i \in NB, \ \rho \sigma_i X_i X_B(r) > 0$$
 where $X_B(r)$ is the r-th column of X_B . Let $L(u)$ denote the indices of $i \in NB$ forming the u smallest ratios. The algorithm then determines the value of u (say \overline{u}) that satisfies

$$|\pi_{r}| - \sum_{i \in L(\overline{u}-1)} \rho \sigma_{i} X_{i} X_{B(r)}^{-1} > 1$$

$$|\pi_{r}| - \sum_{i \in L(\overline{u})} \rho \sigma_{i} X_{i} X_{B(r)} \leq 1.$$

Let t be the unique index which is in $L(\overline{u})$ and not in $L(\overline{u}-1)$. The row of X to enter the basis in the r-th position is X_+ . The algorithm

updates X_B , sets $\sigma_i = -\sigma_i$ for $i \in L(\overline{u}-1)$ and proceeds to the next iteration. The process terminates when $-1 \le \pi_i \le 1$, $i = 1, 2, \ldots, m$.

2.2 Extension of the LAV Algorithm to Handle Dummy Variables

Problem (5) can clearly be handled within the framework of the algorithm outlined in the previous subsection - we need only expand the observation matrix to include the dummy variables. However, we will show that substantial savings can be achieved by recognizing the problem's special structure.

The linear programming equivalent of problem (5) is:

Minimize
$$\sum_{i=1}^{n} (P_i + N_i)$$
 (6)

subject to

$$x_{i1}^{\beta_1} + x_{i2}^{\beta_2} + \dots + x_{im}^{\beta_m} + \sum_{k=1}^{K} d_{ik}^{\alpha_k} + P_i - N_i = y_i;$$
 $P_i \ge 0, N_i \ge 0; i = 1, 2, \dots, n.$

Rewriting the constraints in matrix notation,

$$XB + D\alpha + IP - IN = Y$$
.

Further simplifying the notation, we define

$$V = (X D)$$
 and $\gamma = \begin{pmatrix} \beta \\ \alpha \end{pmatrix}$.

The constraints become

$$V_Y + IP - IN = Y$$
.

The algorithm of Barrodale and Roberts can now be applied directly with V and γ taking the place of X and β , respectively. Let V_B denote a basis for the problem formed by m+K independent nows of V. From the independence of the rows of V_B it follows that at least one V_i with $d_{ik}=1$ (i.e., the m+k-th element of V_i equals 1) must be present in V_B for $k=1, 2, \ldots K$. Hence, V_B may be partitioned as

$$v_B = \begin{pmatrix} x_F & 0_F \\ x_G & 1 \end{pmatrix}$$

where X_F is an m by m submatrix of X, X_G is a K by m submatrix of X, I is a K by K identity matrix and D_F is an m by K matrix with at most one nonzero entry in each row (the nonzero entry being unity).

The inverse of V_B is given by $V_B^{-1} = \begin{bmatrix} (x_F - D_F X_G)^{-1} - (x_F - D_F X_G)^{-1} & D_F \\ -x_G (x_F - D_F X_G)^{-1} & 1 + x_G (x_F - D_F X_G)^{-1} & D_F \end{bmatrix}$

We will demonstrate how the steps of the Barrodale and Roberts' algorithm can be executed conveniently with existing data, $W = \left(X_F - D_F X_G\right)^{-1} \text{ and a knowledge of the rows forming } X_F \text{ and } X_G.$ The logic behind this refined algorithm will be similar in many respects to the generalized upper bounding techniques (see Dantzig and Van Slyke, (1967)) of linear programming. In fact, upon taking the LP dual of (6), generalized upper bound constraints can be recognized.

Proceeding in a manner analogous to that used in subsection 2.1, we define NB to be the index set of nonbisic rows and i(q) to be the q-th basic row. We also define Y_F and Y_G to be the vector with components y_i corresponding to X_F and X_G .

Associated with VR, we have a current solution

$$\binom{\overline{\beta}}{\overline{\alpha}} = V_B^{-1} \binom{Y_F}{Y_G}.$$

Through the previous partitioning scheme,

$$\overline{\beta} = WY_F - WD_FY_G$$

$$\overline{\beta} = W(Y_F - D_FY_G)$$

$$\overline{\beta} = W(Y^*),$$

where the q-th component of Y* equals $y_i(q) = y_i(m+k)$ when $d_i(q), k = 1, k = 1, 2, ..., K$ and equals $y_i(q)$ when $d_i(q), k = 0$, k = 1, 2, ..., K. Now α_k is easily determined by

$$\overline{\alpha}_k = y_{i(m+k)} - X_{i(m+k)}\overline{\beta}.$$

The sign of the residual for the nonbasic rows may be obtained in a similar manner to that given in subsection 2.1. The only difference is an extra addition that is required to evaluate the residual.

The value of the vector of dual variables at any iteration can also be calculated easily from W. We have

$$\pi = (\pi_1, \pi_2, \dots, \pi_m, \pi_{m+1}, \dots, \pi_{m+K}) = (\sum_{i \in NB} \sigma_i(X_i, D_i))V_B^{-1},$$

where D; is the i-th row of D.

By utilizing the partitioning of
$$V_B^{-1}$$
,
 $\pi = (\sum_{i \in NB} \sigma_i(X_i W - D_i X_G W), \sum_{i \in NB} \sigma_i(D_i + D_i W_G W D_F)).$

Let
$$d^* = \sum_{i \in NB} \sigma_i D_i$$
 and $Z = \sum_{i \in NB} \sigma_i X_i - d^* X_G$, then $\pi = (ZW, d^* - ZWD_F)$.

Finally, we define the index set

$$Q(k) = \{q | d_{i(q),k} = 1; k = 1, 2, ..., K; q \le m\}$$

It can now be observed that the first m components of π are equal to ZW and the last K components are given by

$$\pi_{m+k} = d_k^* - \sum_{q \in Q(k)} \pi_q, k = 1, 2, ..., K.$$

The optimality condition is

$$-1 \le \pi_i \le +1$$
, i = 1, 2, ..., m+K.

Assuming the optimality condition is not satisfied, we have $|\pi_n| > 1$ and ρ = sgn (π_n) . Three mutually exclusive and exhaustive cases can arise.

CASE 1. r > m, $d_{i(r), k} = 1$ and $Q(k) = \phi$. This means that the row to leave the basis corresponds to the only observation in the basis affected by factor k. Since the new basis must be nonsingular, the row to enter the basis has the form (X_i, D_i) , $d_{ik} = 1$. The ratio test is simplified to consider

$$|y_i - X_i \overline{\beta} - \overline{\alpha}_i|$$
; $d_{ik} = 1$, ieNB, $\rho \sigma_i > 0$.

Let L(u) denote the indices of i yielding the u smallest of these values. Then $u = [\pi_{\mu}/2]$, where $[\pi_{\mu}/2]$ indicates the greatest integer less than $\pi_{\mu}/2$. Let t be unique index which is in L(u) and not in

L(u-1). The row (X_t, D_t) , enters the basis in the r-th position. Notice that W does not change and the algorithm may proceed directly to the next iteration after updating σ_i , i ε L(u-1).

CASE 2. r < m. The ratios are given by

$$(y_i - X_i \overline{B} - \overline{\alpha}_k)/(\rho(X_i - X_i(m+k))W_r); i \in NB, d_{ik} = 1, k = 1, 2, ..., K$$

 $\rho \sigma_i (X_i - X_i(m+k))W_r > 0$

and

$$(y_i - X_i \overline{\beta})/(\rho X_i W_r)$$
, $i \in NB$; $d_{ik} = 0$, $k = 1, 2, ..., K$; $\rho \sigma_i X_i W_r > 0$ where W_r is the r-th column of W .

The remaining steps are completely analogous to those given previously and will not be repeated.

CASE 3.
$$r > m$$
, $d_{i(r), \overline{k}} = 1$ and $Q(\overline{k}) \neq \phi$

In this situation we reorder the rows of V by interchanging

$$(X_{i(r)}, D_{i(r)})$$

with an $(X_{i(r^*)}, D_{i(r^*)})$, $r^* \in Q(\overline{k})$ (in other words, $d_{i(r^*), \overline{k}} = 1$).

The new W denoted by W* is given by

$$W_{aj}^* = W_{aj}$$
 $q \neq r^*, j = 1, 2, ..., m$

and

$$W_{r*j}^* = \sum_{q \in Q(\overline{k})} W_{qj}, \quad j = 1, 2, \ldots, m.$$

We are now in CASE 2 with $r = r^*$ and may proceed as previously indicated.

3. IMPLEMENTATION AND COMPUTATIONAL RESULTS

Both algorithms described in section 2 have been coded in FORTRAN by the authors. They are maintained as independent subroutines with all input and output as parameters of the CALL statements. We refer to the subroutine corresponding to the revised simplex implementation of Barrodale and Roberts as LINORM. The subroutine corresponding to the algorithm outlined in the subsection 2.2 is called LIDUM. The remainder of this section compares these subroutines with regard to computer and solution times.

3.1 Computer Storage.

Subroutine L1NORM saved the V matrix explicitly. This could easily be reduced with extra coding and without knowledge of the compact inverse procedures discussed here. The same can not be said of the basis inverse. This was represented as an (m+K) by (m+K) matrix in L1NORM and as an m by m matrix in L1DUM.

L1DUM maintained a capacity of K \leq 15, n \leq 300 and m \leq 15. While L1NORM maintained a capacity of m+K \leq 25 and n \leq 300. Some additional coding was, of course, required for L1DUM. Overall L1NORM utilized approximately 1500 more words of internal storage than L1DUM.

We did not employ double precision accuracy on any variables. On some machines (we ran on a CDC 6600 with a sixty bit word) or when solving less stable problems, double precision might be recommended and then the difference in storage utilization would be more dramatic. It should be noted that although we have developed the algorithms using an explicit representation of the inverse, the algorithm can easily be adapted to other methods for solving linear systems. In particular, decomposition methods such as those discussed by Bartels and Golub (1969) and Cline (1976) may be recommended in certain instances.

3.2 Solution Time.

A battery of test problems were solved with two computer codes LINORM and LIDUM. LIDUM and LINORM began with the same initial basis and used the rule of the max $\{|\pi_i|\}$ to define the row to leave the basis at any iteration.

Table I presents a summary of the computational results. Problem sets were randomly generated with five problems in each set. All runs were on the CDC 6600 at The University of Texas at Austin. Times are CPU times in seconds. The calls to the system clock were made in the main program immediately before and immediately after the respective subroutines were called. The number of pivots indicates the number of times the basis inverse was updated.

TABLE I

	PROBLEM SIZE		L1	LDUM	LIN	L1NORM			
#	K	m	n	Time	Pivots	Time	Pivots		
1	4	4	60	.103	13.8	.119	15.0		
2	8	4	60	.105	13.4	.186	17.6		
3 4	12	4	60	.119	15.8	.263	20.2		
	4	8	60	.210	20.4	.216	20.4		
5	8	8	60	.212	22.4	.304	24.2		
6	12	8	60	.213	23.6	.429	26.6		
7	4	4	120	.210	13.0	.235	15.2		
8	8	4	120	.314	20.4	.447	25.4		
9	12	4	120	.298	19.2	.598	28.0		
10	4	8	120	.549	28.6	.506	29.8		
11	8	8	120	.578	32.8	.709	35.4		
12	12	8	120	.566	31.6	.898	36.0		
13	4	4	180	.450	19.2	.491	22.6		
14	8	4	180	.550	22.4	.766	29.4		
15	12	4	180	.476	20.6	.948	33.2		
16	4	8	180	.894	31.2	.826	32.0		
17	8	8	180	1.052	37.8	1.139	39.2		
18	12	8	180	1.128	42.0	1.634	49.8		
19	4	4	240	.613	29.2	.609	20.2		
20	8	4	240	.885	26.2	1.227	35.2		
21	12	4	240	.890	30.4	1.399	38.6		
22	4	8	240	1.326	35.2	1.230	37.2		
23	8	8	240	1.492	40.2	1.608	43.0		
24	12	8	240	1.758	49.0	2.291	54.4		
25	4	4	300	.995	24.4	.998	25.6		
26	8	4	300	1.083	26.4	1.408	33.6		
27	8	4	300	1.125	27.6	1.952	42.6		
28	4	8	300	1.997	43.2	1.768	43.4		
29	8	8	300	2.104	45.4	2.256	49.4		
30	12	8	300	2.423	53.6	3.191	62.4		

This table presents computational results with thirty sets of test problems. Five problems were solved for each dimension and reported times and iterations are means of the results. All times are in CPU seconds.

L1DUM performed better in most instances. As expected, the relative efficiency of L1DUM increased as the number of parameters associated with dummy variables (i.e., the value of K) increased.

4. APPLICATIONS

In problems that require dummy variables, the effect of outliers becomes more difficult to assess as the number of independent variables and dummy variables increases. The motivation for using a robust procedure such as LAV estimation is the same for problems that require dummy variables as it is in the general regression situation. The main purpose of this paper is the refinement of the linear programming approaches to LAV estimation. In this section, we discuss two applications of LAV estimation that involve dummy variables, and a small numerical example is presented.

There are several ways that dummy variable techniques are used in statistical data analysis. One situation occurs when a designed experiment is carried out, and the experimenter is aware of one or more "concomitant variables" that may affect the response and cannot be controlled by the experimenter. This situation is usually referred to as the analysis of covariance and combines features of regression analysis and analysis of variance.

To illustrate the estimation procedure, we use the data in Table II. The response is the average daily gain of pigs, and the experimenter is interested in the effect of four feeds. There are two covariates - initial age and weight - which are believed to affect weight gain. The values of the two covariates and the dummy variables are given in Table II. The LAV estimates of the parameters are given in Table III, which also shows the least squares estimates. If the experimenter wants to evaluate the relative importance of one of the covariates, he can delete that column and recompute the LAV estimates. It is possible to develop various approaches - see e.g., McNeil and Tukey (1975) - to assessing the "qu lity of fit" associated with a given model that can be used in a manner that parallels

TABLE II SOURCE: Snedecor and Cochran (1967, p. 440)

Days	Pounds	Pounds/Day	Feed 1	Feed 2	Feed 3	Feed 4
× _{i1}	×i2	yi	d _{i1}	d ₁₂	d _{i3}	d _{i4}
78 90 94	61 59 76	1.40 1.79 1.72		0 0 0	0 0 0 0	0 0 0
71 99	50 61	1.47 1.26		0	0	0
80 83 75 62 67	54 57 45 41 40	1.28 1.34 1.55 1.57 1.26	1 1 1 1	0 0 0 0	0 0 0 0	0 0 0 0
78 99 80 75 94	74 75 64 48 62	1.61 1.31 1.12 1.35 1.29	0 0 0 0	1 1 1 1	0 0 0 0	0 0 0 0
91 75 63 62 67	42 52 43 50 40	1.24 1.29 1.43 1.29 1.26	0 0 0 0	1 1 1 1	0 0 0 0	0 0 0 0
78 83 79 70 85	80 61 62 47 59	2.67 1.41 1.73 1.23 1.49	0 0 0 0	0 0 0 0	1 1 1 1	0 0 0 0
83 71 66 67 67	42 47 42 40 40	1.22 1.39 1.39 1.46 1.36	0 0 0 0	0 0 0 0	1 1 1 1	0 0 0 0
77 71 78 70 95	62 55 62 43 57	1.40 1.47 1.37 1.15 1.22	0 0 0 0	0 0 0 0	0 0 0 0	1 1 1 1
96 71 63 62 67	51 41 40 45 39	1.48 1.31 1.27 1.22 1.36	0 0 0 0	0 0 0 0	0 0 0 0	1 1 1 1

LAV Estimation For Covariance Model

significance testing in the LS analysis. There is, however, no sampling theory for L1NORM estimates. Consequently, these model-building techniques are most useful in exploratory data analysis.

TABLE III

Parameter

	β ₁	β2	^a 1	α ₂	α ₃	α4
LAV	006706	.009594	1.466	1.326	1.310	1.310
LS	003454	.007414	1.337	1.182	1.318	1.217

Estimates of the parameters for data in Table II.

A second situation in which dummy variables can be used is when adjusting for the effect of seasonal factors in an economic time series. Suppose that the trend-cycle portion of the time series is represented with an "empirical function" composed of polynomial pieces called cubic splines, and that y_i is the observed value of the time series in month i, for $i=1,\ldots,n$. Then the $x_{i,j}$'s in equation (2) are defined as follows:

$$x_{i,j} = i^{j}, (j = 1, 2, 3),$$

$$x_{ij} = \begin{cases} (i-t_{j-3}^*)^3 & \text{if } i \geq t_{j-3}^*, (j = 4, ..., m+3), \\ 0 & \text{otherwise} \end{cases}$$

$$x_{ij} = \begin{cases} 1 & \text{if } i - K \left[(i-1)/K \right] = j - (m+3), (j = m+4, ..., K+m+4) \\ 0 & \text{otherwise,} \end{cases}$$

where the knots $\{t_j^*; j=1,\ldots,m\}$ are known constants which divide the time domain into m+1 intervals, and there are k=12 levels of the seasonal factor. Frome and Armstrong (1977) considered LAV estimation for this model and have presented a numerical example with m=9 and n=132 where the response is the residential construction authorized in Texas for each month in the time period from 1966 to 1976.

5. CONCLUSIONS

With the development of more efficient computer codes (algorithms) to obtain the L_1 norm estimates, the L_1 norm has become an important tool in data analysis. This paper has demonstrated how the Barrodale and Roberts' algorithm can be specialized to solve a class of problems involving dummy variables which arise in several instances. The specialized algorithm has been shown to reduce solution times and computer storage requirements when several columns of the model are associated with dummy variables. Although we have not studied the numerical properties in detail, numerically stable procedures for solving linear systems can easily be combined with the framework of the algorithm.

The two FORTRAN subroutines used to obtain the results reported here are available from the authors for a handling charge.

BIBLIOGRAPHY

Armstrong, R. D. & Frome, E. L. (1976). A comparison of two algorithms for absolute deviation curve fitting. J. Amer. Statist. Assoc. 71, 328-30.

Armstrong, R. D. & Hultz, J. W. (1977). A restricted discrete approximation problem in the L_1 norm. <u>SIAM J. Numer. Analysis 14</u>, 555-65.

Barrodale, I. (1968). L_1 Approximation and the analysis of data. Appl. Statist 17, 51-7.

Barrodale, I. & Roberts, F. D. K. (1973). An improved algorithm for discrete L₁ linear approximation. <u>SIAM J. Numer. Analysis</u> 10, 839-48.

Barrodale, I. & Roberts, F. D. K. (1974). Solution of an over-determined system of equations in the L_1 norm. Commun. Assoc. Comput. Mach. 17, 319-20.

Bartels, R. H. & Golub, G. H. (1969). The simplex method of linear programming using LU decomposition. <u>Commun. Assoc. Comput. Mach.</u> 12, 266-8.

Charnes, A. (1952). Optimality and degeneracy in linear programming. Econometrica 20, 160-70.

Charnes, A. & Cooper, W.W. (1961). Management Models and Industrial Applications of Linear Programming, Vol. I and II. New York: John Wiley and Sons, Inc.

Charnes, A., Cooper, W.W., & Ferguson, R. (1955). Optimal estimation of execution compensation by linear programming. Management Sci. 2, 138-51.

Cline, A.K. (1976). A descent method for the uniform solution to overdetermined systems of linear equations. SIAM J. Numer. Analysis 13, 293-309.

Dantzig, G.B. & Van Slyke, R.M. (1967). Generalized upper bounding techniques. J. Computer Systems Sci. 1, 213-26.

Frome, E.L. & Armstrong, R.D. (1977). A robust procedure for estimating the trend-cycle component of an economic time series. Proc. 10th Ann. Symp. Comput. Sci. Statist. Interface (in press).

Glahe, F.R. & Hunt, J.G. (1970). The small sample properties of simultaneous equation least absolute estimators vis-a-vis least squares estimators. Econometrica 38, 742-53.

Graybill, F.A. (1961). An Introduction to Linear Statistical Models. New York: McGraw-Hill.

Kiountouzis, E.A. (1973). Linear programming techniques in regression analysis. Appl. Statist. 22, 69-73.

McNeil, D.R. & Tukey, J.W. (1975). Higher order diagnosis of two-way tables, illustrated on two sets of demographic empirical distributions. Biometrics 31, 487-510.

Searle, S.R. (1971). <u>Linear Models</u>. New York: John Wiley and Sons, Inc.

Snedecor, G.W. & Cochran, W.G. (1967). Statistical Methods. Iowa State University Press.

Spyropoulos, K., Kiountouzis, E. & Young, A. (1973). Discrete approximations in the L_1 norm. The Computer J. 16, 180-86.

Wagner, H.M. (1959). Linear programming techniques for regression analysis. J. Amer. Statist. Assoc. 54, 206-12.

Security Classification			
DOCUMENT CONTR	ROL DATA - R	& D	
(Security classification of title, body of abstract and indexing a	nnotation must be		
OHIGINATING ACTIVITY (Corporate author)			CURITY CLASSIFICATION
Center for Cybernetic Studies		Unclas	ssified
The University of Texas		26. GROUP	
A Special Purpose Linear Programming	Algorithm	for Chtair	ing Loost Absolute
Value Estimators in a Linear Model with			ing Least Absolute
DESCRIPTIVE NOTES (Type of report and inclusive dates)			acian angan
5. AUTHORISI (First name, middle Initial, last nume)			
Ronald D. Armstrong Edward Frome			Asoga Pardicingo. [Kankingari] basel.]
September 1977	70. TOTAL NO. 0	F PAGES	76. NO. OF REFS
84. CONTRACT OR GRANT NO.	98. ORIGINATOR	S REPORT NUM	BER(\$)
NSF Grant MSC707-00100; N00014-75- 6. PROJECT NO. C-0569		r Cyberne Report Co	tic Studies CS 318
c.	9b. OTHER REPO	RT NO(\$) (Any o	ther numbers that may be easigned
d.			
10. DISTRIBUTION STATEMENT			
This document has been approved for pulits distribution is unlimited.	blic release	and sale;	
II. SUPPLEMENTARY NOTES	12. SPONSURING	MILITARY ACTI	VITY
	Cffice of	Naval Res	earch (Code 434)
	Washingto		
13. ABSTRACT			
A			
Dummy (0, 1) variables are fr	equently us	ed in stati	stical modeling
to represent the effect of certain extrane			
special purpose linear programming algo			
value estimators in a linear model with d			
employs a compact basis inverse procedu	unning vari	ables. In	the advanced basis
exchange techniques available in speciali	and algorith	broorates	the advanced basis
least-absolute-value problem. Computat	nonal resul	ts with a c	omputer code
version of the algorithm are given.			

Unclassified Security Classification

KIY WORDS	THE CAME	LINKA		LINKB		LINKC	
	e certain sin	HOLL	w T	HOLE	wt	HOLL	"
Least Absolute Values							
						Sect 10	
L ₁ Norm							
-1							
Analysis of Covariance							
San división y customas a							
Regression				10.00			3 19
Generalized Upper Bounding							
		HOTELS!					
Linear Programming							
		100			Y- THE		
			0.211				
		150,000					
	DESCRIPTION OF	Time of					
		155 15 3			200		
				114 8 4			
great because an exercispated beingsto		La Visio					
			1000	Carl All			
			e mili	To the last			185

DD . FORM .. 1473 (BACK)

Unclassified
Security Classification